

AD-A192 228

ASYMPTOTICALLY EFFICIENT ADAPTIVE ALLOCATION SCHEMES  
FOR CONTROLLED MARKO. (U) MICHIGAN UNIV ANN ARBOR  
COMMUNICATIONS AND SIGNAL PROCESSING L.

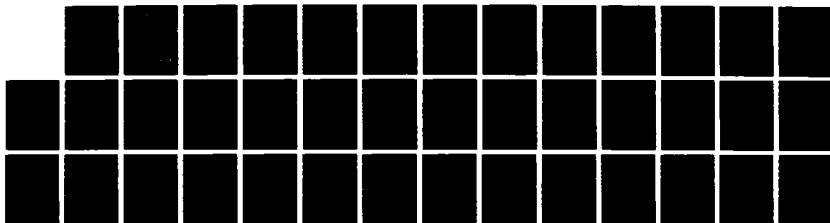
1/1

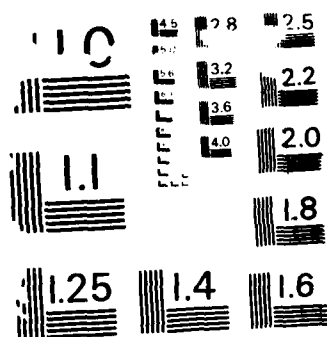
UNCLASSIFIED

R AGRANAL ET AL. FEB 88 TR-254

F/G 12/3

NL





RESOLUTION TEST CHART  
NATIONAL BUREAU OF STANDARDS - 1963-A

AD-A192 228

4

# **ASYMPTOTICALLY EFFICIENT ADAPTIVE ALLOCATION SCHEMES FOR CONTROLLED MARKOV CHAINS: FINITE PARAMETER SPACE**

R. Agrawal, D. Teneketzis

**COMMUNICATIONS & SIGNAL PROCESSING LABORATORY  
Department of Electrical Engineering and Computer Science  
The University of Michigan  
Ann Arbor, MI 48109**

and

V. Anantharam

**School of Electrical Engineering  
Cornell University  
Ithaca, NY 14853**

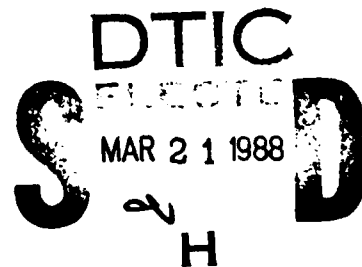
February 1988

Technical Report No. 254  
Approved for public release; distribution unlimited.

Prepared for  
**OFFICE OF NAVAL RESEARCH**  
Department of the Navy  
Arlington, Virginia 22217

and

**NATIONAL SCIENCE FOUNDATION**  
Washington, D.C. 20550



88 3 10 023

## REPORT DOCUMENTATION PAGE

1a. REPORT SECURITY CLASSIFICATION UNCLASSIFIED			1b. RESTRICTIVE MARKINGS NONE		
2a. SECURITY CLASSIFICATION AUTHORITY			3. DISTRIBUTION/AVAILABILITY OF REPORT Approved for Public Release; Distribution Unlimited		
2b. DECLASSIFICATION/DOWNGRADING SCHEDULE					
4. PERFORMING ORGANIZATION REPORT NUMBER(S)  TR 254			5. MONITORING ORGANIZATION REPORT NUMBER(S)		
6a. NAME OF PERFORMING ORGANIZATION Communications & Signal Processing Laboratory		6b. OFFICE SYMBOL (if applicable)	7a. NAME OF MONITORING ORGANIZATION Office of Naval Research and National Science Foundation		
6c. ADDRESS (City, State, and ZIP Code) The University of Michigan Ann Arbor, Michigan 48109-2122			7b. ADDRESS (City, State, and ZIP Code) ONR, 800 N. Quincy St., Arlington, VA 22217 NSF, Washington, D.C. 20550		
8a. NAME OF FUNDING/SPONSORING ORGANIZATION		8b. OFFICE SYMBOL (if applicable)	9. PROCUREMENT INSTRUMENT IDENTIFICATION NUMBER ONR Contract No. N00014-87-K-0540 NSF Grant No. ECS-8517708		
8c. ADDRESS (City, State, and ZIP Code)			10. SOURCE OF FUNDING NUMBERS		
		PROGRAM ELEMENT NO.	PROJECT NO.	TASK NO.	WORK UNIT ACCESSION NO.
11. TITLE (Include Security Classification) Asymptotically Efficient Adaptive Allocation Schemes for Controlled Markov Chains: Finite Parameter Space.					
12. PERSONAL AUTHOR(S) Rajeev Agrawal, Demosthenis Teneketzis, Venkatachalam Anantharam					
13a. TYPE OF REPORT Tech. Report		13b. TIME COVERED FROM TO		14. DATE OF REPORT (Year, Month, Day) February 1988	
15. PAGE COUNT 36					
16. SUPPLEMENTARY NOTATION  → Line ds					
17. COSATI CODES			18. SUBJECT TERMS (Continue on reverse if necessary and identify by block number)		
FIELD	GROUP	SUB-GROUP	Adaptive control scheme.		
			Controlled Markov chain		
			Asymptotically optimal.		
19. ABSTRACT (Continue on reverse if necessary and identify by block number) We consider a controlled Markov chain whose transition probabilities and initial distribution are parametrized by an unknown parameter $\theta$ belonging to some known parameter space $\Theta$ . There is a one-step reward associated with each pair of control and the following state of the process. The objective is to maximize the expected value of the sum of one step rewards over an infinite horizon. By introducing the Loss associated with a control scheme, we show that our problem is equivalent to minimizing this Loss. We define uniformly good adaptive control schemes and restrict attention to these schemes. We develop a lower bound on the Loss associated with any uniformly good control scheme. Finally, we construct an adaptive control scheme, whose Loss equals the lower bound, and is therefore optimal.					
20. DISTRIBUTION/AVAILABILITY OF ABSTRACT <input checked="" type="checkbox"/> UNCLASSIFIED/DUNLIMITED <input type="checkbox"/> SAME AS RPT. <input type="checkbox"/> DTIC USERS			21. ABSTRACT SECURITY CLASSIFICATION UNCLASSIFIED		
22a. NAME OF RESPONSIBLE INDIVIDUAL Carol S. Van Aken			22b. TELEPHONE (Include Area Code) (313) 764-5220		22c. OFFICE SYMBOL

SECURITY CLASSIFICATION OF THIS PAGE

SECURITY CLASSIFICATION OF THIS PAGE

# Asymptotically Efficient Adaptive Allocation Schemes for Controlled Markov Chains: Finite Parameter Space

Rajeev Agrawal, Demosthenis Teneketzis  
Department of Electrical Engineering and Computer Science  
and  
Communications and Signal Processing Laboratory  
University of Michigan  
Ann Arbor, MI 48109-2122

and

Venkatachalam Anantharam  
School of Electrical Engineering  
Cornell University  
Ithaca, NY 14853

February, 1988  
Technical Report No. 254

## Abstract

We consider a controlled Markov chain whose transition probabilities and initial distribution are parametrized by an unknown parameter  $\theta$  belonging to some known parameter space  $\Theta$ . There is a one-step reward associated with each pair of control and the following state of the process. The objective is to maximize the expected value of the sum of one step rewards over an infinite horizon. By introducing the *Loss* associated with a control scheme, we show that our problem is equivalent to minimizing this *Loss*. We define *uniformly good* adaptive control schemes and restrict attention to these schemes. We develop a lower bound on the *Loss* associated with any *uniformly good* control scheme. Finally, we construct an adaptive control scheme whose *Loss* equals the lower bound, and is therefore optimal.

# 1. Introduction

Consider the following stochastic adaptive control problem: The system is modelled by a controlled Markov chain with an unknown parameter, i.e.

$$\mathcal{P}_\theta\{X_{n+1} = y | X_n = x, X_{n-1}, \dots, X_0, U_n, \dots, U_0\} = P(x, y; U_n, \theta) \quad (1.1)$$

where  $X_0, U_0, X_1, U_1, \dots, X_n, U_n, X_{n+1}, \dots$  is the chronological sequence of states and control actions, and  $\theta$  is an unknown parameter belonging to some known parameter space  $\Theta$ ; and

$$\mathcal{P}_\theta(X_0 = x) = p(x; \theta) \quad (1.2)$$

where  $\theta$  is the same as in (1.1). There is a one-step reward  $r(X_n, U_n)$ , associated with each pair  $(X_n, U_n), n \geq 0$ . The objective is to find an adaptive control scheme which maximizes, in some sense, the expected value of the sum of one-step rewards

$$E_\theta J_n = E_\theta \sum_{i=0}^{n-1} r(X_i, U_i), \text{ as } n \rightarrow \infty. \quad (1.3)$$

One of the current approaches to stochastic adaptive control problems is the so called "Certainty Equivalent Control with Forcing" (cf [1]). This scheme is self-tuning in the Cesaro sense and is therefore also optimal for an average reward per unit time criterion (cf [1]). The reward criterion described by (1.3) suggests that we need to determine the maximum rate of increase of  $E_\theta J_n$  as  $n \rightarrow \infty$ . This requirement introduces a notion of optimality that is stronger than the one suggested by the average reward per unit time criterion used in [1] - [7]. For the



Station For	
GFA&I	<input checked="" type="checkbox"/>
TAB	<input type="checkbox"/>
anced	<input type="checkbox"/>
Location	

Distribution/	
Availability Codes	
Avail and/or	
Dist	Special
A-1	

criterion (1.3) it is no longer clear that the Certainty Equivalent Control with Forcing is optimal.

The same reward criterion as (1.3) was previously used in [8] for the study of the controlled i.i.d. process problem. This criterion was initially used by Lai and Robbins [9], [10] for the multi-armed bandit problem. Various extensions of the Lai and Robbins formulation of the multi-armed bandit problems have been reported in [11] and [12]. In this paper we show that the adaptive control problem of Markov chains can be viewed as bandit problem with Markovian rewards. Such a relation provides a convenient way of analyzing the problem, and allows us to develop an "efficient" adaptive control scheme. (We shall precisely define what we mean by efficient in Section 3.)

## 2. The Problem

### 2.1 The System Model

Consider a stochastic system described by a controlled Markov chain on the state space  $\mathcal{X}$ , with control set  $\mathcal{U}$ , transition probability matrix

$$P(u, \theta) := \{P(x, y; u, \theta) | x, y \in \mathcal{X}\} \quad (2.1)$$

and initial probability mass function

$$p(\theta) := \{p(x; \theta) | x \in \mathcal{X}\} . \quad (2.2)$$

The parameter  $\theta$  is unknown, but belongs to a known set  $\Theta$ . Assume that  $\mathcal{X}, \mathcal{U}$



and  $\Theta$  are all finite. Further assume that for

$$x, y \in \mathcal{X}; u \in \mathcal{U}; \theta, \theta' \in \Theta, P(x, y; u, \theta) > 0 \Rightarrow P(x, y; u, \theta') > 0 ; \quad (2.3)$$

for every stationary control law  $g : \mathcal{X} \rightarrow \mathcal{U}$

$$P^g(\theta) := \{P(x, y; g(x), \theta) | x, y \in \mathcal{X}\} \quad (2.4)$$

is irreducible and aperiodic for all  $\theta \in \Theta$  , and

$$p(x; \theta) > 0 \text{ for all } x \in \mathcal{X} \text{ and } \theta \in \Theta . \quad (2.5)$$

Let

$$\pi^g(\theta) := \{\pi^g(x; \theta) | x \in \mathcal{X}\} \quad (2.6)$$

be the stationary distribution corresponding to  $P^g(\theta)$  and let

$$\mu^g(\theta) := \sum_{x \in \mathcal{X}} \pi^g(x; \theta) r(x, g(x)) \quad (2.7)$$

be the mean reward under that stationary distribution.

An "adaptive control scheme"  $\gamma$  is a sequence of random variables  $\{U_n\}_{n=0}^{\infty}$  taking values in the set  $\mathcal{U}$  such that the event  $\{U_n = u\}$  belongs to the  $\sigma$ -field  $\mathcal{F}_n$  generated by  $X_0, U_0, X_1, U_1, \dots, U_{n-1}, X_n$ . Let  $r(X_i, U_i)$  represent the one step reward at time  $i$ , where  $r : \mathcal{X} \times \mathcal{U} \rightarrow R$ . Further define  $J_n := \sum_{i=0}^{n-1} r(X_i, U_i)$  the total reward at time  $n$  as the sum of the one-step rewards upto time  $n$ .

Our objective is to find an adaptive control scheme  $\gamma$  which maximizes, in some sense,  $E_\theta J_n$  as  $n \rightarrow \infty$ . We shall clarify this notion of optimality in Section 2.4. To

achieve our objective we would like to express approximately  $E_\theta J_n$  in terms of the expected number of times each of the stationary control laws  $g$  is used up to time  $n$ , and the expected one-step reward under the invariant distribution corresponding to each  $g$ . For this purpose we need to translate any adaptive control scheme  $\gamma$  to an equivalent adaptive control scheme  $\gamma'$  with the following features:

- (F1) The control scheme  $\gamma'$  chooses a stationary control law  $g_n$  (instead of a control action  $U_n$ ) at each time  $n$ .
- (F2) Whenever a fixed but arbitrary stationary control law  $g$ , chosen by  $\gamma'$ , is used, the sequence of states observed is Markovian. Moreover the sequence of states corresponding to the different stationary control laws, chosen by  $\gamma'$ , are independent conditioned on the initial state.

In Section 2.2 we identify a set of conditions which if satisfied, lead to a control scheme  $\gamma'$  that has the above features, and we construct such an equivalent control scheme. In section 2.3 we define the probability space  $(\Omega', \mathcal{F}', P'_\theta)$  which allows us to define a sequence of states which for each stationary law  $g$  is Markovian, and independent of the sequence of states of any other stationary law  $g'$ , conditioned on all their initial states. Using  $(\Omega', \mathcal{F}', P'_\theta)$  and  $\gamma'$  we can define a control problem that is equivalent to the original one, and we can express  $E_\theta J_n$  in terms of the expected number of times each of the stationary control laws  $g$  is used up to  $n$ , and the expected one-step reward under the invariant distribution corresponding to each  $g$ . Such an expression for  $E_\theta J_n$  allows us to precisely define the sense in

which we want to maximize it.

## 2.2 The Translation Scheme

**Lemma 2.1** Given a controlled Markov chain on a finite state space  $\mathcal{X}$  and with a finite control set  $\mathcal{U}$ , for any adaptive control scheme  $\gamma$  (as defined earlier) there exists an “equivalent adaptive control scheme”  $\gamma'$  taking values on the set  $\mathcal{G} := \{g : \mathcal{X} \rightarrow \mathcal{U}\}$  of stationary control laws with the following properties.

- (i)  $\gamma'$  is a sequence of random variables  $\{g_n\}_{n=0}^{\infty}$  taking values on the set  $\mathcal{G}$  such that the event  $\{g_n = g\}$  belongs to the  $\sigma$ -field  $\mathcal{F}'_n$  generated by  $X_0, g_0, X_1, g_1, \dots, g_{n-1}, X_n$ .
- (ii)  $U_n(\omega) = g_n(X_n)(\omega) \quad \forall n, \omega$ .
- (iii) If  $n_k$  and  $n_{k+1}$  are any two successive time instants at which a stationary control law  $g$  (fixed, but arbitrary) is used, i.e.  $g_{n_k} = g_{n_{k+1}} = g$  and  $g_n \neq g, n_k < n < n_{k+1}$  then  $X_{n_k+1} = X_{n_{k+1}}$ .

(Notice that (i) implies  $\mathcal{F}_n = \mathcal{F}'_n$ .)

### Proof (by construction)

Let  $\#\mathcal{X} = k$  and let  $x^1, x^2, \dots, x^k$  be a prior (but arbitrary) ordering of  $x$ . Similarly let  $\#\mathcal{U} = l$  and  $\mathcal{U} = \{u^1, u^2, \dots, u^l\}$ . To start off observe  $X_0$  and then reorder  $\mathcal{X}$  as  $x^1, x^2, \dots, x^k$  by a left cyclic shift of the prior ordering, such that  $x^1 = X_0$ . Define  $\mathcal{G}_0^i; i = 1, \dots, k$  inductively as follows:

$$\mathcal{G}_0^1 = \{g \in \mathcal{G} : g(x^j) = u^1, 1 < j \leq k\}$$

$$\mathcal{G}_0^i = \{g \in \mathcal{G} : g(x^j) = u^1, i < j \leq k\} - \bigcup_{j=1}^{i-1} \mathcal{G}_0^j; \quad i = 2, \dots, k.$$

Notice that  $\mathcal{G}_0^i; i = 1, \dots, k$  defines a partition of  $\mathcal{G}$ , i.e.  $\bigcup_{i=1}^k \mathcal{G}_0^i = \mathcal{G}$  and  $i \neq j \Rightarrow \mathcal{G}_0^i \cap \mathcal{G}_0^j = \emptyset$ .

Now suppose at time  $n \geq 0$ , i.e. after observing  $X_n$ , we have a partition  $\mathcal{G}_n^i; i = 1, \dots, k$  of  $\mathcal{G}$  with the following five properties:

P1)  $\mathcal{G}_n^i; i = 1, \dots, k$  is determined by  $\mathcal{F}_n'$

P2)  $\forall 1 \leq i \leq k \quad \forall g \in \mathcal{G}_n^i$ , the last time upto time  $n-1$  that the control  $g$  was used (if any) was followed by the state  $x^i$ .

Let

$$X_n = x^{j_n} \quad \text{for some } j_n = 1, \dots, k \quad (2.8)$$

Then,

P3)  $\forall j_n \leq m \leq k$  and for any  $f_m : \{x^1, \dots, x^m\} \rightarrow \mathcal{U}$  there exists a unique  $g \in \bigcup_{i=1}^m \mathcal{G}_n^i \ni g|_{\{x^1, \dots, x^m\}} = f_m$

P4)  $\forall 1 \leq m < j_n$  there exists a unique  $f'_m : \{x^1, \dots, x^m\} \rightarrow \mathcal{U} \ni \forall g \in \bigcup_{i=1}^m \mathcal{G}_n^i, g|_{\{x^1, \dots, x^m\}} \neq f'_m$ , and

P5)  $\forall 1 < m < j_n$  the above found  $f'_m$ 's satisfy  $f'_{m-1} = f'_m|_{\{x^1, \dots, x^{m-1}\}}$

Also assume that

P6)  $g_j$ ,  $0 \leq j < n$  satisfy properties (i), (ii) and (iii) of Lemma 2.1.

We shall now show that we can choose a  $g_n$  satisfying property (P6) on the basis of  $\mathcal{F}'_n$  and construct a new partition  $\mathcal{G}_{n+1}^i; i = 1, \dots, k$  satisfying properties (P1) - (P5) assumed true for time  $n$ . Choose  $g_n \in \mathcal{G}_n^{j_n}$  ( $j_n$  as determined by (2.8)) such that

$$g_n|_{\{x^1, \dots, x^{j_n-1}\}} = f'_{j_n-1} \quad \text{and} \quad g_n(x^{j_n}) = g_n(X_n) = U_n. \quad (2.9)$$

Such a choice is clearly possible by the above induction hypothesis (properties (P3) & (P4)). By noting the fact that  $U_n$  is determined by  $\mathcal{F}_n = \mathcal{F}'_n$  and by the induction hypothesis (properties (P1) and (P2) and (P6)) it follows that (P6) is satisfied for  $n+1$ . Next, let  $X_{n+1} = x^{j_{n+1}}$  for some  $j_{n+1} = 1, \dots, k$ . If  $j_{n+1} = j_n$  then  $\mathcal{G}_{n+1}^i := \mathcal{G}_n^i \quad \forall i = 1, \dots, k$ , and it trivially follows that  $\mathcal{G}_{n+1}^i$ ,  $i = 1, \dots, k$  satisfy (P1)-(P5). Else, if  $j_{n+1} \neq j_n$ ,  $\mathcal{G}_{n+1}^{j_n} := \mathcal{G}_n^{j_n} - \{g_n\}$ ,  $\mathcal{G}_{n+1}^{j_{n+1}} := \mathcal{G}_n^{j_{n+1}} + \{g_n\}$ , and  $\forall i \neq j_n, j_{n+1}$ ,  $\mathcal{G}_{n+1}^i := \mathcal{G}_n^i$ . In this case also it is easy to check that  $\mathcal{G}_{n+1}^i$  satisfy (P1) & (P2). To show that  $\mathcal{G}_{n+1}^i$  satisfy (P3)-(P5) consider two cases

**Case 1**  $j_{n+1} > j_n$  :

-  $\forall j_{n+1} \leq m \leq k \quad \bigcup_{i=1}^m \mathcal{G}_{n+1}^i = \bigcup_{i=1}^m \mathcal{G}_n^i - \{g_n\} + \{g_n\} = \bigcup_{i=1}^m \mathcal{G}_n^i$  Thus (P3) is satisfied.

-  $\forall 1 \leq m < j_n \quad \bigcup_{i=1}^m \mathcal{G}_{n+1}^i = \bigcup_{i=1}^m \mathcal{G}_n^i$ . Thus (P4) & (P5) are satisfied for  $1 \leq m < j_n$  and  $1 < m < j_n$  respectively

-  $\forall j_n \leq m < j_{n+1} \bigcup_{i=1}^m \mathcal{G}_{n+1}^m = \bigcup_{i=1}^m \mathcal{G}_n^m - \{g_n\}$ . Consider the  $f'_m = g_n|_{\{x^1, \dots, x^m\}}$ .

By the induction hypothesis (P3) it then follows that (P4) is satisfied for

$$j_n \leq m < j_{n+1}.$$

Clearly this construction of  $f'_m$  also satisfies

$$f'_{m-1} = f'_m|_{\{x^1, \dots, x^{m-1}\}} \quad \forall j_n < m < j_{n+1}$$

and by (2.9) it also follows that

$$\underset{\text{(old)}}{f'_{j_n-1}} = \underset{\text{(new)}}{f'_{j_n}}|_{\{x^1, \dots, x^{j_n-1}\}}.$$

Thus (P5) is satisfied for  $j_n \leq m < j_{n+1}$ .

**Case 2.**  $j_{n+1} < j_n$

-  $\forall j_n \leq m \leq k \bigcup_{i=1}^m \mathcal{G}_{n+1}^i = \bigcup_{i=1}^m \mathcal{G}_n^i - \{g_n\} + \{g_n\} = \bigcup_{i=1}^m \mathcal{G}_n^i$ . Thus (P3) is satisfied for  $j_n \leq m \leq k$

-  $\forall j_{n+1} \leq m < j_n \bigcup_{i=1}^m \mathcal{G}_{n+1}^i = \bigcup_{i=1}^m \mathcal{G}_n^i + \{g_n\}$ . And since  $f'_m = g_n|_{\{x^1, \dots, x^m\}}$  was the unique one missing from  $\bigcup_{i=1}^m \mathcal{G}_n^i$  (by (2.9) and induction hypothesis (P4), (P5)) it follows that (P3) is now satisfied for  $j_{n+1} \leq m < j_n$ .

-  $\forall 1 \leq m < j_{n+1} \bigcup_{i=1}^m \mathcal{G}_{n+1}^i = \bigcup_{i=1}^m \mathcal{G}_n^i$  and thus (P4) & (P5) are satisfied.

The proof of Lemma 2.1 is now complete (using induction) by checking that the induction hypothesis is satisfied at  $n = 0$ .

## 2.3 Extending the Probability Space

Let  $\Omega = (\mathcal{X} \times \mathcal{U})^\infty$  be the space of all  $\mathcal{X} \times \mathcal{U}$  sequences (i.e. sequences of the type  $X_0, U_0, X_1, U_1, \dots$ ). Give  $(\mathcal{X} \times \mathcal{U})^\infty$  the product  $\sigma$ -field  $\mathcal{F} = \sigma((\mathcal{X} \times \mathcal{U})^\infty)$ , namely, the smallest  $\sigma$ -field such that  $X_0, U_0, X_1, U_1, \dots$  are measurable. There is a unique probability  $\mathcal{P}_\theta^\gamma$  on  $(\Omega, \mathcal{F})$  such that for all  $n$  and all  $x_0, \dots, x_n$  in  $\mathcal{X}$  and  $u_0, \dots, u_n$  in  $\mathcal{U}$ ,

$$\begin{aligned} \mathcal{P}_\theta^\gamma\{X_i = x_i, U_i = u_i, \text{ for } i = 0, 1, \dots, n\} \\ = p(x_0; \theta) \prod_{i=0}^{n-1} P(x_i, x_{i+1}; u_i, \theta) \\ \times \prod_{i=0}^n 1\{\gamma_i(x_0, u_0, \dots, x_i) = u_i\} . \end{aligned} \quad (2.10)$$

This triple  $(\Omega, \mathcal{F}, \mathcal{P}_\theta^\gamma)$  is the minimal underlying probability space required for the description of the problem we address in this paper.

For purposes of analysis and to capture feature (F2) it is useful to extend this probability space which we shall now proceed to do as follows: Let  $\mathcal{G} = \{g^1, \dots, g^d\}$ , and  $\mathcal{X}^d = \{\mathbf{x} = (x^{g^1}, \dots, x^{g^d}) : x^{g^i} \in \mathcal{X}\}$ . Let  $\Omega' = (\mathcal{X}^d)^\infty$  be the space of all  $\mathcal{X}^d$  sequences (i.e. sequences of the type  $\mathbf{X}_0, \mathbf{X}_1, \dots$ ). Give  $(\mathcal{X}^d)^\infty$  the product  $\sigma$ -field  $\mathcal{F}' = \sigma((\mathcal{X}^d)^\infty)$ , namely, the smallest  $\sigma$ -field such that  $\mathbf{X}_0, \mathbf{X}_1, \dots$  are measurable. There is a unique probability  $\mathcal{P}'_\theta$  on  $(\Omega', \mathcal{F}')$  such that for all  $n$  and all  $\mathbf{x}_0, \mathbf{x}_1, \dots, \mathbf{x}_n$  in  $\mathcal{X}^d$ ,

$$\mathcal{P}'_\theta\{X_i = x_i \text{ for } i = 0, 1, \dots, n\}$$

$$= p'_\theta(f(x_0)) \prod_{j=1}^d \prod_{i=0}^{n-1} P^{g^j}(x_i^{g^j}, x_{i+1}^{g^j}; \theta) \quad (2.11)$$

where  $f : \mathcal{X}^d \rightarrow \mathcal{X} \cup \{\Delta\}$ ,  $\Delta$  is an arbitrary element used to augment the state space  $\mathcal{X}$  for the purposes of analysis, and  $f$  is defined as follows: For each  $x \in \mathcal{X}$  left cyclically shift  $\{x^1 \dots x^k\}$  to  $\{x^1, \dots, x^k\}$  such that  $x^1 = x$ . Consider  $\mathcal{G}_0^i$  (from section 2.2) constructed as before on the ordering  $\{x^1, \dots, x^k\}$ . Let  $h : \mathcal{X} \rightarrow \mathcal{X}^d$  such that if  $g^j \in \mathcal{G}_0^i$  then  $h^j(x) = x^i$ . Clearly,  $h$  is one-to-one, but not onto. Let  $h[\mathcal{X}]$  be the range of  $h$ , and  $h^{-1} : h[\mathcal{X}] \rightarrow \mathcal{X}$  be the inverse of  $h$  on its range ( $h^{-1}$  is well-defined as  $h$  is one-to-one.) Finally, let  $f|_{h[\mathcal{X}]} = h^{-1}$  and  $f(x) = \Delta \quad \forall x \in \mathcal{X}^d - h[\mathcal{X}]$ , and  $p'_\theta|_{\mathcal{X}} = p(\theta)$  (defined by (2.2)) and  $p'_\theta(\Delta) = 0$ .

Now on this probability space that we have constructed (note that there is no dependence on the adaptive control scheme  $\gamma$  so far) we can define the random process  $X_0^\gamma, U_0^\gamma, X_1^\gamma, U_1^\gamma, \dots$  by using the equivalent adaptive control scheme  $\gamma'$ . To start off let  $X_0^\gamma := f(X_0)$ . Now given  $X_0^\gamma, U_0^\gamma, \dots, X_n^\gamma$  choose adaptively  $g_n$  such that,  $U_n^\gamma := g_n(X_n^\gamma)$  and  $X_{n+1}^\gamma := X_{T_{g_n}^{g_n}+1}^{g_n}$  where  $T_{g_n}^{g_n}$  is the number of times the control law  $g_n$  was used upto time  $n$  (in  $X_0, U_0, \dots, X_n$ ), and  $X_{T_{g_n}^{g_n}+1}^{g_n}$  is the component of  $X_{T_{g_n}^{g_n}+1}$  corresponding to  $g_n$ . It can be easily verified that the random process  $X_0^\gamma, U_0^\gamma, X_1^\gamma, U_1^\gamma, \dots$  constructed above has the same distribution (in  $(\Omega', \mathcal{F}', \mathcal{P}'_\theta)$ ) as the one given by  $(\Omega, \mathcal{F}, \mathcal{P}_\theta^\gamma)$ . Note that for  $X_0 \ni f(X_0) = \Delta$  the process is undefined, but that is not important as  $\mathcal{P}'_\theta\{X_0 : f(X_0) = \Delta\} = 0$ .

Using  $(\Omega', \mathcal{F}', \mathcal{P}'_\theta)$  and  $\gamma'$  we can now express  $E_\theta J_n$  in terms of the expected



number of times each stationary control law  $g$  is used and the expected one-step reward under the invariant distribution corresponding to each  $g$ .

## 2.4 Analysis of the Reward Criterion

Consider

$$\begin{aligned}
 J_n &= \sum_{i=0}^{n-1} r(X_i, U_i) \\
 &= \sum_{i=0}^{n-1} r(X_i, U_i) \sum_{g \in \mathcal{G}} 1(g_i = g) \sum_{x \in \mathcal{X}} 1(X_i = x) \\
 &= \sum_{g \in \mathcal{G}} \sum_{x \in \mathcal{X}} \sum_{i=0}^{n-1} r(X_i, U_i) 1(g_i = g) 1(X_i = x) \\
 &= \sum_{g \in \mathcal{G}} \sum_{x \in \mathcal{X}} r(x, g(x)) N^g(x, T_n^g)
 \end{aligned} \tag{2.12}$$

where

$$\begin{aligned}
 N^g(x, T_n^g) &= \sum_{i=0}^{T_n^g-1} 1(X_i^g = x) \\
 &= \sum_{i=0}^{n-1} 1(X_i = x, g_i = g)
 \end{aligned}$$

and

$$T_n^g = \sum_{i=0}^{n-1} 1(g_i = g) . \tag{2.13}$$

Note that in the extended probability space  $(\Omega', \mathcal{F}', \mathcal{P}_\theta')$   $T_n^g$  is a stopping w.r.t. the increasing family of  $\sigma$ -algebras  $\{(\bigvee_{\substack{g' \in \mathcal{G} \\ g' \neq g}} \mathcal{F}_\infty^{g'}) \bigvee \mathcal{F}_n^g\}$  where  $\mathcal{F}_n^g = \sigma(X_0^g, X_1^g, \dots, X_n^g)$

and  $\mathcal{F}_\infty^g = \bigvee_n \mathcal{F}_n^g$ .

To express  $EN^g(x, T_n^g)$  in terms of the invariant distribution under  $g$  and  $ET_n^g$  we use the following result:

**Lemma 2.2** Let  $X_0, X_1, \dots$  be Markovian with finite state space  $\mathcal{X}$ , transition matrix  $P$ -irreducible and aperiodic and stationary distribution  $\pi$ . Let  $\mathcal{F}_n$  denote the  $\sigma$ -algebra generated by  $X_0, X_1, \dots, X_n$ . Let  $\mathcal{G}$  be another  $\sigma$ -algebra and  $A$  an event such that  $A \in \mathcal{F}_0 \vee \mathcal{G}$  and  $\{X_0 = x\} \cap A = \begin{cases} A & A \subset \{X_0 = x\} \\ \phi & \text{otherwise} \end{cases}$ . Furthermore let  $\mathcal{G}$  be independent of  $\mathcal{F}_\infty$  conditioned on the event  $A$ . Let  $\tau$  be a stopping of  $\{\mathcal{G} \vee \mathcal{F}_n\}$  such that  $E[\tau|A] < \infty$ . Let

$$N(x, \tau) = \sum_{i=0}^{\tau-1} 1(X_i = x)$$

Then, for some fixed constant  $K$ , independent of  $A, x$  and  $\tau$ .

$$|E[N(x, \tau)|A] - \pi(x)E[\tau|A]| \leq K \quad (2.14)$$

**Proof:** Follows from Lemma 2.1 in [11].

Notice that  $\bigvee_{\substack{g' \in \mathcal{G} \\ g' \neq g}} \mathcal{F}_\infty^{g'}$  and  $\mathcal{F}_\infty^g$  are independent conditioned on the event  $A_x = \{X_0 = x\}, x \in \mathcal{X}^d$ . Moreover  $A_x \in \bigvee_{g \in \mathcal{G}} \mathcal{F}_0^g \subset ((\bigvee_{\substack{g' \in \mathcal{G} \\ g' \neq g}} \mathcal{F}_\infty^{g'}) \vee \mathcal{F}_0^g)$  and  $\{X_0^g = x\} \cap \{X_0 = x\} = \begin{cases} \{X_0 = x\}; \{X_0 = x\} \subset \{X_0^g = x\}. \\ \phi \quad \text{otherwise} \end{cases}$ .

Therefore by Lemma 2.2 it follows that

$$|E_\theta[N^g(x, T_n^g)|A_x] - \pi^g(x; \theta)E_\theta[T_n^g|A_x]| \leq K$$

for some fixed constant  $K$  independent of  $\underline{x}$ ,  $x$  and  $n$ .

Thus,

$$|E_\theta[N^g(x, T_n^g)] - \pi^g(x, \theta)E_\theta[T_n^g]| \leq K \quad (2.15)$$

From (2.13) and (2.15) it follows that

$$|E_\theta J_n - \sum_{g \in \mathcal{G}} \mu^g(\theta) E_\theta T_n^g| \leq K' \quad (2.16)$$

where  $K'$  is independent of  $n$  and  $\mu^g(\theta)$  is as defined by (2.7). Let  $g^*(\theta) = \arg \max_{g \in \mathcal{G}} (\mu^g(\theta))$ , and for simplicity assume that it is unique for each  $\theta \in \Theta$ . Thus if we knew the true parameter the control scheme  $g_n = g^*(\theta)$  gives the optimal reward (upto a constant) for all  $n$ , and for this scheme

$$|E_\theta J_n - n\mu^{g^*(\theta)}(\theta)| \leq K'.$$

In the absence of the knowledge of the true parameter it is desirable to approach this performance as closely as possible. For this purpose we define the *Loss* associated with an adaptive control scheme  $\gamma$ ,

$$L_n(\theta) := n\mu^{g^*(\theta)}(\theta) - E_\theta J_n \quad (2.17)$$

By (2.16) it follows that

$$|L_n(\theta) - \sum_{\substack{g \in \mathcal{G} \\ g \neq g^*(\theta)}} (\mu^{g^*(\theta)}(\theta) - \mu^g(\theta)) E_\theta T_n^g| \leq \text{const.} \quad (2.18)$$

Maximizing  $E_\theta J_n$  is thus equivalent to minimizing the *Loss*. More precisely we want to minimize the rate at which the *Loss* increases with  $n$  (e.g. finite, logarithmic,

linear etc.). Thus, this is a stronger criterion for optimality than the average reward per unit time criterion (used in [1] - [7]) which only requires the *Loss* to be  $o(n)$ . In view of (2.18) the above problem is reduced to one of minimizing the rate at which  $E_\theta T_n^g$  increases for  $g \in \mathcal{G}, g \neq g^*(\theta)$ .

Note that it is impossible to minimize  $L_n(\theta)$  uniformly over all parameters  $\theta \in \Theta$ . For example the stationary control scheme  $g_n = g^*(\theta)$  for all  $n$ , will have a finite *Loss* where the true parameter is  $\theta$ . However, when the true parameter is  $\theta'$  such that  $g^*(\theta') \neq g^*(\theta)$ , then this scheme will have a *Loss* proportional to  $n$ . Having made this observation we call a scheme "uniformly good" if for every parameter  $\theta \in \Theta$

$$L_n(\theta) = o(n^\alpha) \text{ for every } \alpha > 0 \quad (2.19)$$

Such schemes do not allow the *Loss* to increase very rapidly for any  $\theta \in \Theta$ . We restrict our attention to the class of uniformly good schemes and consider any others as uninteresting.

### 3. A Lower Bound on the *Loss*

In this section we obtain a lower bound on the *Loss*  $L_n(\theta)$  for certain values of the parameter  $\theta \in \Theta$ . Before we present the bound we introduce the necessary notation. Let

$$B(\theta) := \{\theta' \in \Theta : P^{g^*(\theta)}(\theta') = P^{g^*(\theta)}(\theta) \text{ and } g^*(\theta') \neq g^*(\theta)\} ,$$

$$\mathcal{G}_\theta := \mathcal{G} - \{g^*(\theta)\} ,$$

$$\begin{aligned}
\mathcal{A}_\theta &:= \left\{ (\alpha^g, g \in \mathcal{G}_\theta) : \alpha^g \geq 0, \sum_{g \in \mathcal{G}_\theta} \alpha^g = 1 \right\}, \\
d_\theta(g) &:= (\mu^{g^*(\theta)}(\theta) - \mu^g(\theta)) \text{ and} \\
I^g(\theta, \theta') &:= \sum_{x \in \mathcal{X}} \pi^g(x; \theta) \sum_{y \in \mathcal{X}} P^g(x, y; \theta) \log \frac{P^g(x, y; \theta)}{P^g(x, y; \theta')}.
\end{aligned} \tag{3.1}$$

Note that  $I^g(\theta, \theta')$  is just the expectation with respect to the invariant measure of  $P^g(\theta)$  of the Kulback Liebler numbers between the individual rows of  $P^g(\theta)$  and  $P^g(\theta')$  thought of as probability distributions on  $\mathcal{X}$ .

The bound is now presented in the form of Theorem 3.1 below.

**Theorem 3.1** Let  $\theta \in \Theta$  be such that  $B(\theta)$  is non-empty. Then for any uniformly good control scheme  $\phi$ , under the parameter  $\theta$ ,

$$1) \quad \lim_{n \rightarrow \infty} P_\theta \left\{ \sum_{g_\theta} T_n^g d_\theta(g) < \frac{\log n}{1 + 2\rho} \cdot \frac{1}{\max_{\alpha \in \mathcal{A}_\theta} \min_{\theta' \in B(\theta)} \frac{\sum_{g_\theta} \alpha^g I^g(\theta, \theta')}{\sum_{g_\theta} \alpha^g d_\theta(g)}} \right\} = 0 \quad \forall \rho > 0. \tag{3.2}$$

Consequently,

$$2) \quad \liminf_{n \rightarrow \infty} \frac{L_n(\theta)}{\log n} \geq \min_{\alpha \in \mathcal{A}_\theta} \max_{\theta' \in B(\theta)} \frac{\sum_{g_\theta} \alpha^g d_\theta(g)}{\sum_{g_\theta} \alpha^g I^g(\theta, \theta')}. \tag{3.3}$$

### Proof

The proof can easily be obtained from that of Theorem 3.1 of [8] by substituting  $g$  for  $u$  and  $\mathcal{G}_\theta$  for  $\mathcal{U}_\theta$  and by invoking the ergodic theorem instead of the strong law of large numbers.  $\square$

Note that we do not have a lower bound for those values of  $\theta$  for which  $B(\theta)$  is empty. In view of this observation and the above lower bound we call a scheme “efficient” if

$$\begin{aligned} \limsup_{n \rightarrow \infty} \frac{L_n(\theta)}{\log n} &\leq \min_{\alpha \in \mathcal{A}_\theta} \max_{\theta' \in B(\theta)} \frac{\sum_{g \in \mathcal{G}_\theta} \alpha^g d_\theta(g)}{\sum_{g \in \mathcal{G}_\theta} \alpha^g I^g(\theta, \theta')} \text{ if } B(\theta) \text{ is non-empty} \\ L_n(\theta) &< \infty \text{ if } B(\theta) \text{ is empty} \end{aligned} \quad (3.4)$$

## 4. The Control Scheme

### 4.1 Preliminaries

Let  $M^{(2)}$  be the unit simplex in  $\mathbf{R}^{|\mathcal{X}|^2}$  identified with the space of probability measures on  $\mathcal{X}^2$ .

Let

$$\nu_\theta^g(x, y) := \pi^g(x; \theta) P^g(x, y; \theta); \quad x, y \in \mathcal{X} \quad (4.1)$$

Then  $\nu_\theta^g = \{\nu_\theta^g(x, y) : x, y \in \mathcal{X}\} \in M^{(2)}$ . Since  $\Theta$  and  $\mathcal{G}$  are finite  $\nu_\theta^g$  take on only a finite number of points in  $M^{(2)}$ . Therefore it is possible to find an  $\epsilon > 0$  such for all values of  $\nu_\theta^g$  we can identify  $\epsilon$ -neighborhoods (“ $\epsilon$ -nbd of  $\nu_\theta^g$ ”) of the type:

$$\epsilon\text{-nbd}(\nu_\theta^g) := \{\nu \in M^{(2)} : \max_{x, y \in \mathcal{X}} |\nu(x, y) - \nu_\theta^g(x, y)| < \epsilon\} \quad (4.2)$$

which are disjoint for distinct values of  $\nu_\theta^g$ .

Also define

$$S(\theta) := \{\theta' \in \Theta : P^{g^*(\theta)}(\theta') = P^{g^*(\theta)}(\theta) \text{ and } g^*(\theta') = g^*(\theta)\} \quad (4.3)$$

This is the set of parameters for which the optimal control laws are the *same* as that for  $\theta$ , and the transition probabilities under the optimal control law are also identical. Let

$$\mathcal{G}(S(\theta)) := \{g : P^g(\theta') \neq P^g(\theta), \theta' \in S(\theta)\}. \quad (4.4)$$

Recall from Section 3 that

$$B(\theta) := \{\theta' \in \Theta : P^{g^*(\theta)}(\theta') = P^{g^*(\theta)}(\theta) \text{ and } g^*(\theta') \neq g^*(\theta)\}. \quad (4.5)$$

This is the set of parameters for which the optimal control laws are *better* than the optimal control law for  $\theta$ , and the transition probabilities under the optimal control law for  $\theta$  are identical.

Let

$$\alpha(\theta) = \{\alpha^g(\theta) : g \in \mathcal{G}_\theta\} \quad (4.6)$$

achieve the minimum in the lower bound for the *Loss* in (3.2), where  $\mathcal{G}_\theta = \mathcal{G} - \{g^*(\theta)\}$  and

$$T_{\theta, x_0}^g = E_\theta^g[\inf\{n \geq 1 | X_n = x_0\} | X_0 = x_0], \quad (4.7)$$

be the expected recurrence time of the state  $x_0$  under the control law  $g$ . On the basis of these define,

$$\beta(\theta) = \{\beta^g(\theta) : g \in \mathcal{G}_\theta\} \text{ with } \beta^g(\theta) = \frac{\alpha^g(\theta)/T_{\theta, x_0}^g}{\sum_{g \in \mathcal{G}_\theta} \alpha^g(\theta)/T_{\theta, x_0}^g}. \quad (4.8)$$

## 4.2 Description of the Control Scheme

Let  $x_0 \in \mathcal{X}$  be an arbitrary but fixed state. Define the  $\{\mathcal{F}_t = \sigma(X_0, U_0, X_1, \dots, X_{t-1}, U_{t-1}, X_t)\}$  stopping times  $\tau_0, \tau_1, \dots$  by  $\tau_m := \inf\{t > \tau_{m-1} | X_t = x_0\}$ ,  $m \geq 1$ , and  $\tau_0 = \inf\{t | X_t = x_0\}$ . The control scheme we construct chooses a stationary control law at times  $0, \tau_0, \tau_1, \dots$  adaptively on the basis of all the past observations and past actions, and use this control law till  $\tau_0 - 1, \tau_1 - 1, \tau_2 - 1, \dots$  respectively. That is, over each recurrence interval marked by the state  $x_0$  we use the same control law which is chosen adaptively at the beginning of that block. With this in mind we now describe how the choice of control laws is made at the beginning of each block. From now on we shall refer to the actual time as time and the recurrence points as instances. Initially, i.e. at  $t = 0$ , choose a fixed but arbitrary control law  $g_0$  and use it till time  $\tau_0 - 1$ . Then to start off, use each of the control laws  $g \in \mathcal{G}$  once each. From then at each recurrence point, compute the empirical pair measure  $\rho_n^g := \{\rho_n^g(x, y) | x, y \in \mathcal{X}\} \in M^{(2)}$  corresponding to each  $g \in \mathcal{G}$  as

$$\rho_n^g(x, y) := \frac{1}{T_n^g - \tau_0} \sum_{i=\tau_0}^{n-1} 1\{g_i = g, X_i = x, X_{i+1} = y\} \quad (4.9)$$

where  $n$  is the actual time

Define the conditions

$C1(\theta)$ :  $\rho_n^g \in \epsilon\text{-nbd}(\nu_\theta^g) \ \forall \ g \in \mathcal{G}$  and  $B(\theta)$  is empty

$C2(\theta)$ :  $\rho_n^g \in \epsilon\text{-nbd}(\nu_\theta^g) \ \forall g \in \mathcal{G}$  and  $B(\theta)$  is non-empty.

$C3$ : there does not exist  $\theta \in \Theta$  such that  $\rho_n^g \in \epsilon\text{-nbd}(\nu_\theta^g) \ \forall g \in \mathcal{G}$ .



(Note that  $C3 = (\bigcup_{\theta \in \Theta} (C1(\theta) \cup C2(\theta)))'$ ). Proceed as follows.

- 1) If  $C1(\theta)$  is satisfied for some  $\theta \in \Theta$  then use  $g^*(\theta)$ .
- 2) If  $C2(\theta)$  is satisfied for some  $\theta \in \Theta$  then do the following: Maintain a count of the number of instances condition  $C2(\theta)$  is satisfied. Of these, for the first instance choose among those control laws  $g \in \mathcal{G}_\theta$  randomly with probabilities  $\beta^g(\theta)$ . Refer to this process as "randomization". For those instances when this count is even (call this situation  $C2(\theta)$  a) use  $g^*(\theta)$ . For other instances when the count is odd (call this situation  $C2(\theta)$  b) compute the likelihood ratio

$$\Lambda_n(\theta) := \lambda_{T_n}(\theta) := \min_{\theta' \in B(\theta)} \prod_{i=0}^{T_n-1} \frac{P^{g_i^r}(X_i^r, X_{i+1}^r; \theta)}{P^{g_i^r}(X_i^r, X_{i+1}^r; \theta')}$$

of  $\theta$  vs  $B(\theta)$ , where  $X_0^r, g_0^r, X_1^r, \dots, g_{T_n-1}^r, X_{T_n}^r$  is the sequence of pairs of control laws used and states observed upto time  $n$  when "randomization" is done with  $\beta(\theta)$ . If  $\Lambda_n > K_{n+1}$  (say  $C2(\theta)b1$ ), where  $K_n = n(\log n)^p$  for some fixed  $p > 1$ , the use  $g^*(\theta)$ . If  $\Lambda_n \leq K_{n+1}$  (say  $C2(\theta)b2$ ) then do the following: Maintain a count of the number of instances this condition ( $C2(\theta)b2$ ) is satisfied. If this count is a perfect square (say  $C2(\theta)b2a$ ) then use round robin amongst  $g \in \mathcal{G}(S(\theta))$ . If this count is not a perfect square (say  $C2(\theta)b2b$ ) then do "randomization" using  $\beta(\theta)$ .

- 3) If  $C3$  is satisfied then use round-robin amongst  $g \in \mathcal{G}$ .

### 4.3 Upper Bound on the Loss

In this section we derive an upper bound on the *Loss* associated with the adaptive control scheme  $\gamma^*$  constructed in Section 4.2. The bound is given by the main Theorem 4.2. Lemmas 4.1, 4.2, 4.3 and Theorem 4.1 are needed for the proof of the main theorem.

**Lemma 4.1:** Let  $X_0, X_1, \dots$  be Markovian with finite state space  $\mathcal{X}$ , transition matrix  $P$ , invariant distribution  $\pi$ , and initial distribution  $p$ . Let  $M^{(2)}$  be the unit simplex on  $R^{|\mathcal{X}|^2}$  identified with the space of probability measures on  $\mathcal{X}^2$ , and let  $K \subset M^{(2)}$ , closed, such that  $\pi P \notin K$ . Let  $\rho_n := \{\rho_n(x, y) | x, y \in \mathcal{X}\}$  where  $\rho_n(x, y) := \frac{1}{n} \sum_{i=0}^{n-1} 1\{X_i = x, X_{i+1} = y\}$ . Then

(i)  $P(\rho_n \in K) < Ae^{-an}$  for all  $n \geq 1$  for some positive constants  $A, a$ .

Let  $N := \sum_{n=1}^{\infty} 1(\rho_n \in K)$ . Then

(ii)  $EN < \infty$

Let  $L := \sup\{n \geq 1 | \rho_n \in K\}$ . Then

(iii)  $EL < \infty$

**Proof:**

Part (i) follows from the theory of large deviations. See [14], Problem IX.6.12.

$$\begin{aligned} EN &= \sum_{n=1}^{\infty} P(\rho_n \in K) \\ &\leq \sum_{n=1}^{\infty} Ae^{-an} \end{aligned}$$

$< \infty$  which proves (ii)

$$\begin{aligned}
 EL &= E \sum_{n=1}^{\infty} 1(\exists i \geq n, \rho_i \in K) \\
 &= E \sum_{n=1}^{\infty} 1(\bigcup_{i \geq n} (\rho_i \in K)) \\
 &\leq \sum_{n=1}^{\infty} \sum_{i=n}^{\infty} P(\rho_i \in K) \\
 &\leq \sum_{n=1}^{\infty} \sum_{i=n}^{\infty} A e^{-ai} \\
 &< \infty \quad \text{which proves (iii),} \quad \square
 \end{aligned}$$

**Lemma 4.2:** Let  $S_n = X_1 + \dots + X_n$  where  $X_1, X_2, \dots$  are i.i.d.,  $EX_1 > 0$  and let  $N = \sum_{n=1}^{\infty} 1(S_n \leq 0)$ ,  $L = \sum_{n=1}^{\infty} 1(\inf_{i \geq n} S_i \leq 0)$ . Then the following are equivalent:

(a)  $E(|X_1|^2 1(X_1 \leq 0)) < \infty$ .

(b)  $EN < \infty$ .

(c)  $EL < \infty$ .

**Proof:** See Hogan [15].

**Lemma 4.3:** Let  $X_1, X_2, \dots$  be i.i.d. Let  $f^i$  be a real valued Borel function such that  $0 < Ef^i(X_1) < \infty, i \in I$ , finite. Let  $S_n^i = f^i(X_1) + f^i(X_2) + \dots + f^i(X_n)$ ,  $L_A^i = \sum_{n=1}^{\infty} 1(\inf_{i \geq n} S_i^i \leq A)$ , and  $L_A = \max_{i \in I} L_A^i$ . If  $E(|f^i(X_1)|^2 1(f^i(X_1) \leq 0)) < \infty$  for all  $i \in I$ , then

$$\limsup_{A \rightarrow \infty} \frac{EL_A}{A} \leq \frac{1}{\min_{i \in I} (Ef^i(X_1))} \quad (4.10)$$

**Proof:** For  $\varepsilon > 0$ , and for any fixed  $i \in I$

$$L_A^i \leq \frac{A(1+\varepsilon)}{Ef^i(X_1)} + L^i \quad (4.11)$$

where

$$L^i = \sum_{n=1}^{\infty} 1 \left( \inf_{t \geq n} \left( S_t^i - \frac{tEf^i(X_1)}{1+\varepsilon} \right) \leq 0 \right) \quad (4.12)$$

Consider the i.i.d. r.v.'s

$$Z_t^i = f^i(X_t) - \frac{Ef^i(X_1)}{1+\varepsilon}$$

We have,

$$\begin{aligned} E\{|Z_1^i|^2 1(Z_1^i \leq 0)\} &\leq 2E \left\{ \left[ |f^i(X_1)|^2 + \left( \frac{Ef^i(X_1)}{1+\varepsilon} \right)^2 \right] 1 \left( f^i(X_1) \leq \frac{Ef^i(X_1)}{1+\varepsilon} \right) \right\} \\ &\leq 2E \left\{ |f^i(X_1)|^2 1(f^i(X_1) \leq 0) \right\} \\ &\quad + 2E \left\{ |f^i(X_1)|^2 1 \left( 0 < f^i(X_1) \leq \frac{Ef^i(X_1)}{1+\varepsilon} \right) \right\} + 2 \left( \frac{Ef^i(X_1)}{1+\varepsilon} \right)^2 \\ &< \infty . \end{aligned}$$

Then, by Lemma 4.2 it follows that  $EL^i < \infty$ .

Therefore

$$E \left( \max_{i \in I} L^i \right) \leq E \left( \sum_{i \in I} L^i \right) = \sum_{i \in I} EL^i = k(\varepsilon) < \infty , \quad (4.13)$$

for some constant  $k(\varepsilon)$  independent of  $A$ .

Now,

$$\begin{aligned} L_A &= \max_{i \in I} L_A^i \leq \max_{i \in I} \left( \frac{A(1+\varepsilon)}{Ef^i(X_1)} + L^i \right) \\ &\leq \frac{A(1+\varepsilon)}{\min_{i \in I} (Ef^i(X_1))} + \max_{i \in I} L^i \end{aligned} \quad (4.14)$$

By (4.11) and (4.12) it follows that

$$EL_A \leq \frac{A(1+\varepsilon)}{\min_{i \in I}(Ef^i(X_1))} + k(\varepsilon)$$

$$\limsup_{A \rightarrow \infty} \frac{EL_A}{A} \leq \frac{1+\varepsilon}{\min_{i \in I}(Ef^i(X_1))}$$

By letting  $\varepsilon \rightarrow 0$  we get the desired result.  $\square$

**Theorem 4.1** Let  $\theta \in \Theta$  be such that  $B(\theta)$  is non-empty. Then,

$$(1) \quad \limsup_{n \rightarrow \infty} \left[ E_\theta \left[ \sum_{m=1}^{\infty} 1(\lambda_{r_m}(\theta) \leq K_{n+1}) \right] / \log n \right] \leq \frac{1}{\min_{\theta' \in B(\theta)} \sum_{\theta''} \beta^{\theta''}(\theta) T_{\theta, x_0}^{\theta''} I^{\theta''}(\theta, \theta')}$$

(4.15)

$$(2) \quad P_{\theta'} \{ \lambda_i(\theta) > K_{n+1} \text{ for some } 1 \leq i \leq n \} \leq \frac{1}{K_{n+1}} \quad \text{for } \theta' \in B(\theta).$$

(4.16)

**Proof:**

Let  $X_0^r, X_1^r, \dots$  be the sequence of observed states when "randomization" is used with  $\alpha(\theta)$ . Let  $X^* = \bigcup_{t \geq 1} X^t$ , with the Borel  $\sigma$ -algebra of the discrete topology, i.e. all subsets are measurable. The process  $\{X_t\}_{t \geq 0}$  allows us to define  $X^*$  valued random variables  $B_1, B_2, \dots$  called blocks as follows: Define the  $\{\mathcal{F}_t\}$  stopping times  $\tau_k, k > 1$  by

$$\tau_k = \inf \{ t > \tau_{k-1} | X_t^r = X_0^r = x_0 \}$$

with  $\tau_0 = 0$ . (Note that  $\tau_k \leq \infty$  a.s.). Then

$$B_k = (X_{\tau_{k-1}}^r, X_{\tau_{k-1}+1}^r, \dots, X_{\tau_k}^r)$$

Let  $B'_k = (B_k, g_k)$ . Since the same control law is used over the entire block, and the choice of the specific law for each block is made by independent randomizations at the beginning of the block it can be easily shown that  $\{B'_k\}$  are i.i.d.

Let

$$f^{\theta'}(B'_k) = \log \frac{P^{g_k}(X_{\tau_{k-1}}^r, X_{\tau_{k-1}+1}^r; \theta) \dots P^{g_k}(X_{\tau_k}^r, X_{\tau_k}^r; \theta)}{P^{g_k}(X_{\tau_{k-1}}^r, X_{\tau_{k-1}+1}^r; \theta') \dots P^{g_k}(X_{\tau_k}^r, X_{\tau_k}^r; \theta')}$$

Then

$$\begin{aligned} E_{\theta}[f^{\theta'}(B'_k)|X_0 = x_0] &= \sum_{\theta} \beta^g(\theta) E_{\theta} \left[ \sum_{t=\tau_{k-1}}^{\tau_k-1} \log \frac{P^g(X_t, X_{t+1}; \theta)}{P^g(X_t, X_{t+1}; \theta')} | X_0 = x_0 \right] \\ &= \sum_{\theta} \beta^g(\theta) E_{\theta} \left[ \sum_{x,y \in \mathcal{X}} N(x, y, B_k) \log \frac{P^g(x, y; \theta)}{P^g(x, y; \theta')} | X_0 = x_0 \right] \\ &= \sum_{\theta} \beta^g(\theta) \sum_{x,y \in \mathcal{X}} \pi^g(x; \theta) P^g(x, y; \theta) T_{\theta, x_0}^g \log \frac{P^g(x, y; \theta)}{P^g(x, y; \theta')} \\ &= \sum_{\theta} \beta^g(\theta) T_{\theta, x_0}^g I^g(\theta, \theta') \end{aligned}$$

and

$$\begin{aligned} E_{\theta}[(f^{\theta'}(B'_k))^2 1(f^{\theta'}(B'_k) \leq 0) | X_0 = x_0] \\ &= \sum_{\theta} \beta^g(\theta) E_{\theta}[(f^{\theta'}((B_k, g)))^2 1(f^{\theta'}((B_k, g)) \leq 0) | X_0 = x_0] \\ &= \sum_{\theta} \beta^g(\theta) \sum_{B_k \in X^*} P^g(B_k; \theta | X_0 = x_0) \left( \log \frac{P^g(B_k; \theta | X_0 = x_0)}{P^g(B_k; \theta' | X_0 = x_0)} \right)^2 \end{aligned}$$

$$\begin{aligned}
& \cdot 1 \left( \log \frac{P^g(B_k; \theta | X_0 = x_0)}{P^g(B_k; \theta' | X_0 = x_0)} \leq 0 \right) \\
&= \sum_{g_\theta} \beta^g(\theta) \sum_{B_k \in X^*} P^g(B_k; \theta' | X_0 = x_0) \frac{P^g(B_k; \theta | X_0 = x_0)}{P^g(B_k; \theta' | X_0 = x_0)} \\
&\quad \cdot \left( \log \frac{P^g(B_k; \theta | X_0 = x_0)}{P^g(B_k; \theta' | X_0 = x_0)} \right)^2 1 \left( \frac{P^g(B_k; \theta | X_0 = x_0)}{P^g(B_k; \theta' | X_0 = x_0)} \leq 1 \right) \\
&\leq \sum_{g_\theta} \beta^g(\theta) \sum_{B_k \in X^*} P^g(B_k; \theta' | X_0 = x_0) \frac{4}{e^2} \text{ as } x(\log x)^2 \leq \frac{4}{e^2} \text{ on } 0 \leq x \leq 1 \\
&= \frac{4}{e^2} < \infty
\end{aligned}$$

Thus by Lemma 4.3 we have the desired result (i).

To prove (ii) note that

$$\begin{aligned}
& \{ \Lambda_i(\theta) > k_{n+1} \text{ for some } 1 \leq i \leq n \} \\
&= \left\{ \min_{\theta' \in B(\theta)} \prod_{t=0}^{i-1} \frac{P^{g_t^r}(X_t^r, X_{t+1}^r; \theta)}{P^{g_t^r}(X_t^r, X_{t+1}^r; \theta')} > k_{n+1} \text{ for some } 1 \leq i \leq n \right\} \\
&\subseteq \left\{ \prod_{t=0}^{i-1} \frac{P^{g_t^r}(X_t^r, X_{t+1}^r; \theta)}{P^{g_t^r}(X_t^r, X_{t+1}^r; \theta')} > k_{n+1} \text{ for some } 1 \leq i \leq n \right\} \\
&\quad \text{for any } \theta' \in B(\theta), \text{ and } \left\{ \prod_{t=0}^{i-1} \frac{P^{g_t^r}(X_t^r, X_{t+1}^r; \theta)}{P^{g_t^r}(X_t^r, X_{t+1}^r; \theta')} \right\}_{t \geq 1} \text{ is a } \mathcal{F}_i \text{ martingale}
\end{aligned}$$

under  $\theta'$  with mean 1.

Thus the result follows by the submartingale inequality (see [13], pg 243).  $\square$

**Theorem 4.2:** Under the adaptive control scheme  $\phi^*$ , for  $g \neq g^*(\theta)$

$$\begin{aligned}
\text{(i)} \quad E_\theta T_n^g &\leq \left( \frac{\alpha^g(\theta)}{\min_{\theta' \in B(\theta)} \sum_{g_\theta} \alpha^g(\theta) I^g(\theta, \theta')} + o(1) \right) \log n \text{ if } B(\theta) \text{ is non-empty} \\
E_\theta T_n^g &< \infty \quad \text{if } B(\theta) \text{ is empty.} \quad (4.17)
\end{aligned}$$

Consequently

$$(ii) \quad L_n(\theta) \leq \left( \max_{\theta' \in B(\theta)} \frac{\sum_{\mathcal{G}_\theta} \alpha^g(\theta) d_\theta(g)}{\sum_{\mathcal{G}_\theta} \alpha^g(\theta) I^g(\theta, \theta')} + o(1) \right) \log n \text{ if } B(\theta) \text{ is non-empty}$$

$$L_n(\theta) < \infty \quad \text{if } B(\theta) \text{ is empty} \quad (4.18)$$

where  $\alpha(\theta) = \{\alpha^g(\theta) : g \in \mathcal{G}_\theta\}$  is defined by (4.6).

**Proof:** As in Section 4.2 define the  $\{\mathcal{F}_t (= \sigma(X_0, U_0, X_1, \dots, X_{t-1}, U_{t-1}, X_t))\}$  stopping times  $\tau_0, \tau_1, \dots$  by  $\tau_m := \inf\{t > \tau_{m-1} | X_t = x_0\}$  with  $\tau_0 = \inf\{n | X_n = x_0\}$ .

Then  $\tau_m < \infty$  a.s.. Then for any  $n \geq 0$ , any  $g \in \mathcal{G}_\theta$  we have

$$T_n^g = \sum_{i=0}^{n-1} 1(g_i = g)$$

$$\leq \sum_{i: \tau_i < n} 1(g_{\tau_i} = g)(\tau_{i+1} - \tau_i) + \tau_0$$

since the choice of  $g$ 's is only made at the stopping times  $\tau_i$ . So

$$E_\theta T_n^g \leq E_\theta \sum_{i=0}^{\infty} 1(g_{\tau_i} = g)(\tau_{i+1} - \tau_i) 1(\tau_i < n) + E_\theta \tau_0$$

$$= \sum_{i=0}^{\infty} E_\theta [E_\theta [1(g_{\tau_i} = g) 1(\tau_i < n)(\tau_{i+1} - \tau_i) | \mathcal{F}_{\tau_i}]] + E_\theta \tau_0$$

$$= \sum_{i=0}^{\infty} E_\theta [1(g_{\tau_i} = g) 1(\tau_i < n) E_\theta [(\tau_{i+1} - \tau_i) | \mathcal{F}_{\tau_i}]] + E_\theta \tau_0$$

$$= \sum_{i=0}^{\infty} E_\theta [1(g_{\tau_i} = g) 1(\tau_i < n) T_{\theta, x_0}^g] + E_\theta \tau_0$$

$$= T_{\theta, x_0}^g E_\theta \sum_{i: \tau_i < n} 1(g_{\tau_i} = g) + E_\theta \tau_i.$$

Let us now examine the term  $\sum_{i: \tau_i < n} 1(G_i = g)$ , where  $G_i = g_{\tau_i}$ .



$$\begin{aligned}
& \sum_{i: \tau_i < n} 1(G_i = g) \\
&= 1 + \sum_{i \geq d: \tau_i < n} 1(G_i = g) \\
&= 1 + \sum_{i \geq d: \tau_i < n} 1\{G_i = g, C1(\theta') \text{ is satisfied at stage } i \text{ for some } \theta' \in \Theta\} \\
&\quad + \sum_{i \geq d: \tau_i < n} 1\{G_i = g, C2(\theta') \text{ is satisfied at stage } i \text{ for some } \theta' \in \Theta\} \\
&\quad + \sum_{i \geq d: \tau_i < n} 1\{G_i = g, C3 \text{ is satisfied at stage } i\} \\
&= 1 + \text{Term 1} + \text{Term 2} + \text{Term 3 (say)}, \tag{4.19}
\end{aligned}$$

where  $C1(\theta')$ ,  $C2(\theta')$  and  $C3$  are defined in Section 4.2 and  $d$  is the cardinality of the set  $\mathcal{G}$  of stationary controls. Let us now examine each term separately. Defining  $\mathcal{L}^g$  by

$$\mathcal{L}^g := \sup_{T_n^g \geq 1} \{\rho_n^g \notin \epsilon\text{-nbd}(\nu_\theta^g)\} . \tag{4.20}$$

and noting that  $E_\theta \mathcal{L}^g < \infty$  by Lemma 4.1(ii), we get

Term 3  $\leq \sum_{g \in \mathcal{G}} \mathcal{L}^g$ , thus,

$$E_\theta \text{Term 3} \leq \sum_{g \in \mathcal{G}} E_\theta \mathcal{L}^g < \infty , \tag{4.21}$$

and Term 1  $\leq \mathcal{L}^g$ , thus,

$$E_\theta \text{Term 1} \leq E_\theta \mathcal{L}^g < \infty . \tag{4.22}$$

$$\text{Term 2} = \sum_{i \geq d: \tau_i < n} 1\{G_i = g, C2(\theta') \text{ is satisfied at stage } i \text{ for some}$$

$$\begin{aligned}
& \theta' \in \Theta \text{ such that } \nu_{\theta'}^{g^*(\theta')} \neq \nu_{\theta}^{g^*(\theta')} \} \\
& + \sum_{i \geq d: \tau_i < n} 1\{G_i = g, C2(\theta') \text{ is satisfied at stage } i \text{ for some} \\
& \quad \theta' \in \Theta \text{ such that } \theta \in B(\theta')\} \\
& + \sum_{i \geq d: \tau_i < n} 1\{G_i = g, C2(\theta') \text{ is satisfied at stage } i \text{ for some} \\
& \quad \theta' \in \Theta \text{ such that } \theta \in S(\theta')\} \\
& + \sum_{i \geq d: \tau_i < n} 1\{G_i = g, C2(\theta) \text{ is satisfied at stage } i\} . \\
& = \text{Term 2a} + \text{Term 2b} + \text{Term 2c} + \text{Term 2d} \text{ (say)} . \tag{4.23}
\end{aligned}$$

Next we upper bound each of terms 2a - 2d separately.

$$\begin{aligned}
\text{Term 2a} &= \sum_{\substack{\theta': B(\theta') \text{ is empty and} \\ \nu_{\theta'}^{g^*(\theta')} \neq \nu_{\theta}^{g^*(\theta')}}} \sum_{i \geq d: \tau_i < n} 1\{G_i = g, C2(\theta') \text{ is satisfied at stage } i\} \\
&\leq \sum_{\substack{\theta': B(\theta') \text{ is not empty and} \\ \nu_{\theta'}^{g^*(\theta')} \neq \nu_{\theta}^{g^*(\theta')}}} \left[ 1 + \sum_{i \geq d: \tau_i < n} 1\{G_i = g^*(\theta'), C2(\theta') \text{ is satisfied at stage } i\} \right] \\
&\leq \sum_{\substack{\theta': B(\theta') \text{ is not empty and} \\ \nu_{\theta'}^{g^*(\theta')} \neq \nu_{\theta}^{g^*(\theta')}}} (\mathcal{L}^{g^*(\theta')} + 1) \tag{4.24}
\end{aligned}$$

The first of the inequalities of (4.24) holds because under  $C2(\theta')$ ,  $g^*(\theta')$  is chosen on all the even instances, therefore, on at least as many instances as any other control minus one. The second of the inequalities of (4.24) holds because the sum on the left hand side counts a subset of the times when  $g^*(\theta')$  is used and  $\rho_n(g^*(\theta')) \notin \epsilon\text{-nbd}(\nu_{\theta}^{g^*(\theta')})$  where  $\theta$  is the true parameter.

By Lemma 4.1(ii) it follows that

$$E_\theta \text{ Term 2a} \leq \sum_{\substack{\theta': B(\theta') \text{ is empty and} \\ \nu_{\theta'}^{g^*}(\theta') \neq \nu_\theta^{g^*}(\theta')}} (1 + E_\theta \mathcal{L}^{g^*}(\theta')) < \infty \quad (4.25)$$

$$\begin{aligned} \text{Term 2b} &\leq \sum_{\theta': \theta \in B(\theta')} \sum_{i \geq d: \tau_i < n} \{C2(\theta') \text{ is satisfied at stage } i\} \\ &\leq \sum_{\theta': \theta \in B(\theta')} 2 \left[ 1 + \sum_{i \geq d: \tau_i < n} 1\{C2(\theta')b \text{ is satisfied at stage } i\} \right] \\ &= \sum_{\theta': \theta \in B(\theta')} 2 \left[ 1 + \sum_{i \geq d: \tau_i < n} 1\{C2(\theta')b1 \text{ is satisfied at stage } i\} \right. \\ &\quad \left. + \sum_{i \geq d: \tau_i < n} 1\{C2(\theta')b2 \text{ is satisfied at stage } i\} \right] \\ &\leq \sum_{\theta': \theta \in B(\theta')} 2 \left[ 1 + \sum_{i \geq d: \tau_i < n} 1\{\Lambda_{\tau_i}(\theta') > K_{\tau_i+1}\} \right. \\ &\quad \left. + \sum_{i \geq d: \tau_i < n} 1\{C2(\theta')b2 \text{ is satisfied at stage } i\} \right] \\ &\leq \sum_{\theta': \theta \in B(\theta')} 2 \left[ 1 + \sum_{i=d}^{\infty} 1\{\lambda_j(\theta') > K_i \text{ for some } j \leq i-1\} \right. \\ &\quad \left. + \sum_{i \geq d: \tau_i < n} 1\{C2(\theta')b2 \text{ is satisfied at stage } i\} \right] \quad (4.26) \end{aligned}$$

The first of the inequalities of (4.26) results by removing the condition  $G_i = g$ . The second one results by observing that the total number of time instants that  $C2(\theta')$  is satisfied is upperbounded by twice the odd instants that  $C2(\theta')$  holds, and by noting that the first time we randomize and the other odd times we call  $C2(\theta')b$ . The third inequality results because  $\{C2(\theta')b2 \text{ is satisfied at stage } i\}$  implies  $\{\Lambda_{\tau_i}(\theta') > K_{\tau_i+1}\}$ .

Consider now the term  $\sum_{i \geq d: \tau_i < n} 1\{C2(\theta')b2 \text{ is satisfied at stage } i\}$ .

$$\begin{aligned}
& \sum_{i \geq d: \tau_i < n} 1\{C2(\theta')b2 \text{ is satisfied at stage } i\} \\
&= \sum_{i \geq d: \tau_i < n} 1\{C2(\theta')b2a \text{ is satisfied at stage } i\} + \sum_{i \geq d: \tau_i < n} 1\{C2(\theta')b2b \text{ is satisfied at stage } i\} \\
&\leq 1 + 2 \sum_{i \geq d: \tau_i < n} 1\{C2(\theta')b2b \text{ is satisfied at stage } i\} \\
&= 1 + 2 \sum_{i \geq d: \tau_i < n} 1\{C2(\theta')b2b \text{ is satisfied at stage } i; \text{ of the number of instances} \\
&\quad \text{that } C2(\theta')b2b \text{ has been satisfied so far, the fraction of instances} \\
&\quad \text{that } g' \text{ is chosen } \in (\beta^{g'}(\theta') - \epsilon, \beta^{g'}(\theta') + \epsilon)\} \\
&\quad + 2 \sum_{i \geq d: \tau_i < n} 1\{C2(\theta')b2b \text{ is satisfied at stage } i; \text{ of the number of instances} \\
&\quad \text{that } C2(\theta')b2b \text{ has been satisfied so far, the fraction of instances} \\
&\quad \text{that } g' \text{ is chosen } \notin (\beta^{g'}(\theta') - \epsilon, \beta^{g'}(\theta') + \epsilon)\} \\
&\leq 1 + 2 \sum_{j=1}^{\infty} 1\{\rho_i(g') \notin \epsilon\text{-nbd}(\nu_{\theta'}^{g'}) \text{ for some } i > (\beta^{g'}(\theta') - \epsilon)j\} \\
&\quad + 2 \sum_{j=1}^{\infty} 1\{\text{of } j \text{ the fraction of instances } g' \text{ is chosen } \notin (\beta^{g'}(\theta') - \epsilon, \beta^{g'}(\theta') + \epsilon)\} \quad (4.27)
\end{aligned}$$

where  $g' \in \mathcal{G}_{\theta'}$  is such that  $\nu_{\theta'}^{g'} \neq \nu_{\theta'}^{g'}$ .

The first of the inequalities of (4.27) results by observing that the number of instances when condition  $C2(\theta')b2a$  is satisfied (i.e. the count of the number of instances  $C2(\theta')b2$  is satisfied is a perfect square) is upper bounded by the number of instances when condition  $C2(\theta')b2b$  is satisfied plus one. Consider now changing the index of summation to the instances when randomization is done. Then the condition  $C2(\theta')b2b$  along with the condition that the fraction of instances that  $g'$  is chosen  $\in (\beta^{g'}(\theta') - \epsilon, \beta^{g'}(\theta') + \epsilon)$  at stage  $i$ , imply that  $\rho_i(g') \notin \epsilon\text{-nbd}(\nu_{\theta'}^{g'})$  for some  $i > (\beta^{g'}(\theta') - \epsilon)j$ . By extending the summation to infinity together with the

above observation establishes the last of the inequalities of (4.27).

Thus, by Lemma 4.1(i) and (4.16) it follows that

$$\begin{aligned}
E_{\theta} \text{ Term } 2b &\leq \sum_{\theta': \theta \in B(\theta')} 2 \left[ 1 + \sum_{i=d}^{\infty} (i(\log i)^p)^{-1} + 1 \right. \\
&\quad \left. + 2 \sum_{j=1}^{\infty} \sum_{i > (\beta \theta'(\theta') - \epsilon)j} A_1 e^{-a_1 i} + 2 \sum_{j=1}^{\infty} A_2 e^{-a_2 j} \right] \\
&< \infty \tag{4.28}
\end{aligned}$$

where  $A_1, a_1, A_2, a_2 > 0$  are some constants.

$$\begin{aligned}
\text{Term } 2c &= \sum_{\theta': \theta \in S(\theta')} \sum_{i \geq d: \tau_i < n} \{G_i = g, C2(\theta') \text{ is satisfied at stage } i\} \\
&\leq \sum_{\theta': \theta \in S(\theta')} \left[ 1 + \sum_{i \geq d: \tau_i < n} 1 \{G_i = g, C2(\theta') \text{ is satisfied at stage } i\} \right] \\
&\leq \sum_{\theta': \theta \in S(\theta')} \left[ 1 + \sum_{i \geq d: \tau_i < n} 1 \{C2(\theta') \text{ is satisfied at stage } i\} \right] \\
&\leq \sum_{\theta': \theta \in S(\theta')} \left[ 1 + l^2 + \sum_{j=1}^{\infty} 1 \{ \rho_j(g') \notin \epsilon\text{-nbd}(\nu_{\theta}^{g'}) \} (2j+1)l^2 \right] \tag{4.29}
\end{aligned}$$

where  $g' \in \mathcal{G}(S(\theta'))$  is such that  $\nu_{\theta}^{g'} \neq \nu_{\theta'}^{g'}$  and  $\# \mathcal{G}(S(\theta')) = l$ .

The first inequality of (4.29) results by noting that since  $\theta \in S(\theta')$ ,  $g \neq g^*(\theta') = g^*(\theta)$  can be chosen only when condition  $C2(\theta')b2$  is satisfied, or at the first instance when  $C2(\theta')$  is true. The second inequality results by removing the requirement  $G_i = g$ . The third inequality results by upperbounding the number of instances condition  $C2(\theta')b2$  is satisfied. This can be achieved as follows: First restrict attention to those instances that are perfect squares and the control  $g'$  is used. At these instances since  $C2(\theta')$  is satisfied  $\rho_n(g') \in \epsilon\text{-nbd}(\nu_{\theta}^{g'})$ , thus, by the choice

of  $g' \in \mathcal{G}(S(\theta'))$ ,  $\rho_n(g') \notin \epsilon\text{-nbd}(\nu_{\theta}^{g'})$ . Consider the sum of the intervals between the above instances. (Note that the length of the  $j^{\text{th}}$  interval is upperbounded by  $[(j+1)^2 - j^2]l^2 = (2j+1)l^2$ .) Then the number of instances condition  $C2(\theta')b2$  is satisfied cannot exceed this sum. Finally, the inequality results by changing the summation index to all the times when  $g'$  is used and upperbounding the interval following the time  $\rho_j(g') \notin \epsilon\text{-nbd}(\nu_{\theta}^{g'})$  by  $(2j+1)l^2$ .

Again, by using Lemma 4.1(i) we get

$$E_{\theta} \text{ Term } 2c \leq \sum_{\theta': \theta \in S(\theta')} \left[ 1 + l^2 + \sum_{j=1}^{\infty} A e^{-aj} \cdot (2j+1)l^2 \right] < \infty \quad (4.30)$$

Now if  $B(\theta)$  is empty then,

$$\text{Term } 2d = 0 \quad (4.31)$$

Otherwise,

$$\begin{aligned} \text{Term } 2d &= \sum_{i \geq d: \tau_i < n} 1\{G_i = g, C2(\theta) \text{ is satisfied at stage } i\} \\ &\leq 1 + \sum_{i \geq d: \tau_i < n} 1\{G_i = g, C2(\theta)b2 \text{ is satisfied at stage } i\} \\ &= 1 + \sum_{i \geq d: \tau_i < n} 1\{G_i = g, C2(\theta)b2a \text{ is satisfied at stage } i\} \\ &\quad + \sum_{i \geq d: \tau_i < n} 1\{G_i = g, C2(\theta)b2b \text{ is satisfied at stage } i\} \\ &\leq 2 + \sum_{i \geq d: \tau_i < n} 1\{G_i = g, C2(\theta)b2b \text{ is satisfied at stage } i\} \\ &\quad + \left( \sum_{i \geq d: \tau_i < n} 1\{C2(\theta)b2b \text{ is satisfied at stage } i\} \right)^{1/2} \end{aligned} \quad (4.32)$$

The first of the inequalities of (4.32) is obtained by noting  $g \neq g^*(\theta)$  can be chosen only at the first instance when  $C2(\theta)$  is satisfied (in which case randomization is

done) or when  $C2(\theta)b2$  is satisfied. The last of the inequalities of (4.32) results because the number of instances condition  $C2(\theta)b2a$  is satisfied is upperbounded by one plus the square root of the number of instances  $C2(\theta)b2b$  is satisfied.

To upperbound  $E_\theta$  Term 2d we use (4.32), Jensen's inequality and the following fact: At each instance  $i$  when condition  $C2(\theta)b2b$  is satisfied, the choice of the control law  $G_i \in \mathcal{G}_\theta$  is made by an independent randomization  $\beta(\theta)$ . Then,

$$\begin{aligned} E_\theta \text{ Term } 2d &\leq 2 + \sum_{i \geq d: \tau_i < n} P_\theta \{C2(\theta)b2b \text{ is satisfied at stage } i\} \cdot \beta^g(\theta) \\ &\quad + \left( \sum_{i \geq d: \tau_i < n} P_\theta \{C2(\theta)b2b \text{ is satisfied at stage } i\} \right)^{1/2} \\ &\leq 2 + \beta^g(\theta) E_\theta [\sup \{1 \leq i \leq n | \lambda_i(\theta) \leq K_{n+1}\}] \\ &\quad + (E_\theta [\sup \{1 \leq k \leq n | \lambda_k(\theta) \leq K_{n+1}\}])^{1/2} \end{aligned} \quad (4.33)$$

Using (4.15) we get

$$\limsup_{n \rightarrow \infty} E_\theta \text{ Term } 2d / \log n \leq \frac{\beta^g(\theta)}{\min_{\theta' \in B(\theta)} \sum_{G_\theta} \beta^g(\theta) T_{\theta, x_0}^g I^g(\theta, \theta')} . \quad (4.34)$$

Combining (4.19), (4.21), (4.22), (4.23), (4.25), (4.28), (4.30), (4.31) and (4.34) we get (4.17). (4.18) follows easily from (4.17) and (2.18).

□

## 5. Conclusions

In this paper we considered the problem of adaptive control of Markov Chains. The optimality criterion used, namely minimizing the rate at which the Loss in-

creases is stronger than the average reward per unit time criterion. Multi-armed bandit problems with “Loss” as the optimality criterion is one class of stochastic adaptive control problems that has previously been analyzed. Therefore one way to proceed with our problem is to relate it to the multi-armed bandit problem, like was done in [8] for the controlled i.i.d. process problem. The translation scheme and the extended probability space are crucial in allowing us to view the adaptive control of Markov chains as a multi-armed bandit problem. The stationary control laws correspond to the “arms”, and the sequence of states observed when any particular stationary control law is used are Markovian. The formulation then resembles that of the multi-armed bandit problem in [11], part II. One very important difference between our problem and that of [11] is that the parametrization of the “arms” in our problem is not independent. This difference is reflected in the lower bound on the Loss we obtain in Section 3, and also needs to be kept in mind when designing an optimal scheme like the one of Section 4. The control scheme presented in Section 4 has an intuitively appealing structure as it clearly specifies the conditions under which there is either only identification, or only control, or identification and control, and treats each one of these conditions optimally.

## Acknowledgements

The research of Rajeev Agrawal and Demosthenis Teneketzis was supported in part by NSF Grant No. ECS-8517708 and ONR Grant No. N00014-87-K-0540.



## References

- [1] P.R. Kumar and P. Varaiya, "Stochastic Systems: Estimation, Identification and Adaptive Control", Prentice-Hall, 1986.
- [2] P. Mandl, "Estimation and Control in Markov Chains", *Adv. Appl. Prob.*, 6 (1974), pp. 40-60.
- [3] V. Borkar and P. Varaiya, "Adaptive Control of Markov Chains, I: finite parameter set", *IEEE Transactions on Automatic Control*, AC-24, 1979, pp. 953-958.
- [4] V. Borkar and P. Varaiya, "Identification and Adaptive Control of Markov Chains", *SIAM J. on Control and Optimization*, 20, 1982, pp. 470-489.
- [5] P. R. Kumar and A. Becker, "A new family of optimal adaptive controllers for Markov chains", *IEEE Transactions on Automatic Control*, AC-27, 1982, pp. 137-146.
- [6] P. R. Kumar and W. Lin, "Optimal adaptive controllers for unknown Markov chains", *IEEE Transactions on Automatic Control*, AC-27, 1982, pp. 765-774.
- [7] R. A. Miloto and J. B. Cruz, "An optimization oriented approach to the adaptive control of Markov chains", *IEEE Transactions on Automatic Control*, Vol AC-32, No. 9, September 1987.
- [8] R. Agrawal, D. Teneketzis, V. Anantharam, "Asymptotically Efficient Allocation Schemes for Controlled I.I.D. Processes: Finite Parameter Space," Technical Report No. 253, Communications and Signal Processing Lab, Univ. of Michigan, January, 1988.
- [9] T.L. Lai and H. Robbins, "Asymptotically Efficient Adaptive Allocation Rules", in *Advances in Applied Mathematics*, 1984.
- [10] T.L. Lai and H. Robbins, "Asymptotically Optimal Allocation of Treatments in Sequential Experiments," In 'Design of Experiments' (eds. T.J. Santner and A.C. Tamhane), Marcel-Dekker, New York.
- [11] V. Anantharam, P. Varaiya, and J. Walrand, "Asymptotically Efficient Allocation Rules for the Multiarmed Bandit Problem with Multiple Plays; Part I: IID Rewards, Part II: Markovian Rewards", *IEEE Transaction on Automatic Control*, Vol 1 AC-32, No. 11, November 1987, pp 968-982.
- [12] R. Agrawal, M. Hegde, and D. Teneketzis, "Asymptotically Efficient Adaptive Allocation Rules for the Multi-armed Bandit Problem with Switching Cost", Technical Report No. 246, Communications and Signal Processing Lab, Univ. of Michigan, April, 1987.
- [13] S. Ross, "Stochastic Processes", Wiley, 1983.

- [14] R.S. Ellis, "Entropy, Large Deviations, and Statistical Mechanics", Springer-Verlag, 1985.
- [15] M. Hogan, "Moments of the Minimum of a Random Walk and Complete Convergence," Technical Report No. 21, Department of Statistics, Stanford University, Jan 1983.

END

DATE

FILMED

6-1988

DTIC